

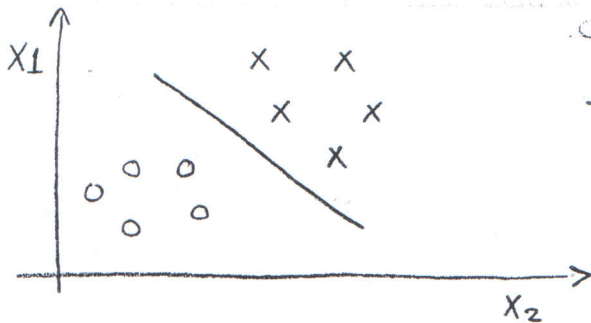
Lec : 8

* UnSupervised Learning

For more information go to page 11

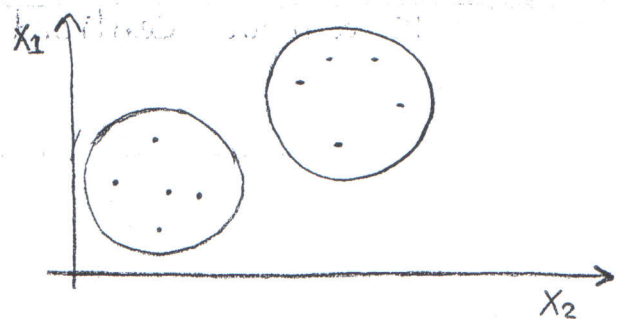
1. given unlabeled training set

Supervised Learning



Given Labeled training set
 $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
↓
right answer

Unsupervised Learning



Given unlabeled training set
 $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$
there is no right answer.

→ One of the most Famous unsupervised learning algorithms is Clustering Algorithm.

Clustering Algorithm. group unlabeled data into different clusters. (according to the relation between the data or determined thing)

Applications OF Clustering

Market Segmentation

Social network

Organize Computing

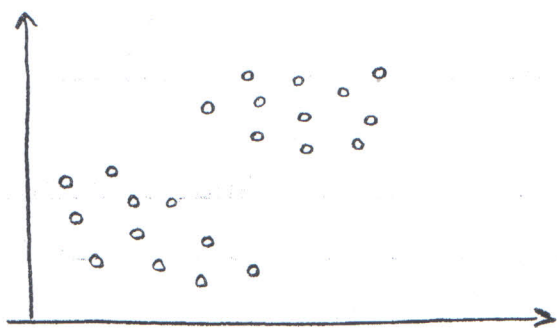
Astronomical data analysis

Clustering

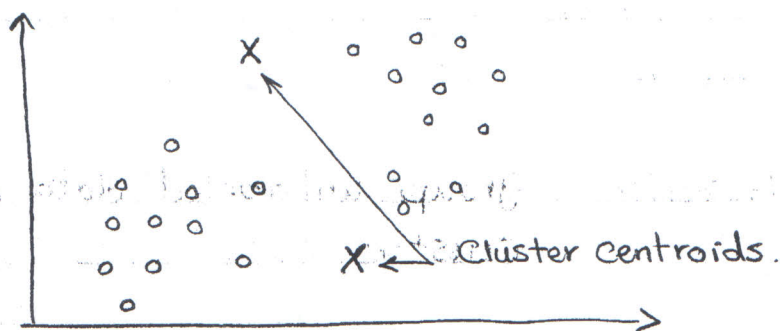
- K-means Algorithm (mathematical algorithm) *
group the unlabeled data into different clusters
According to the distance between each data
Location.

→ K Means is an iterative algorithm and it does two things. First is a cluster assignment step, and second is a move Centroid step.

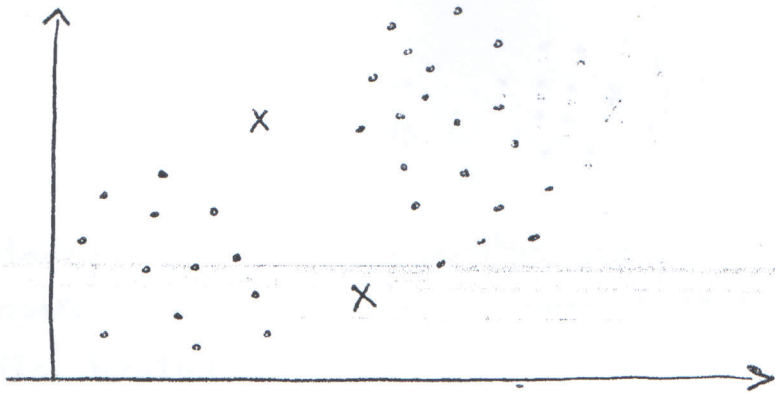
→ We have unlabeled data



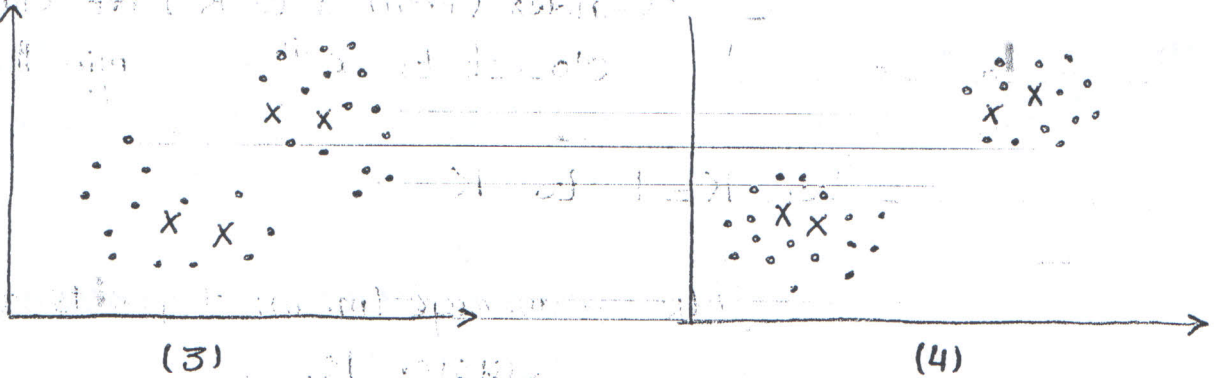
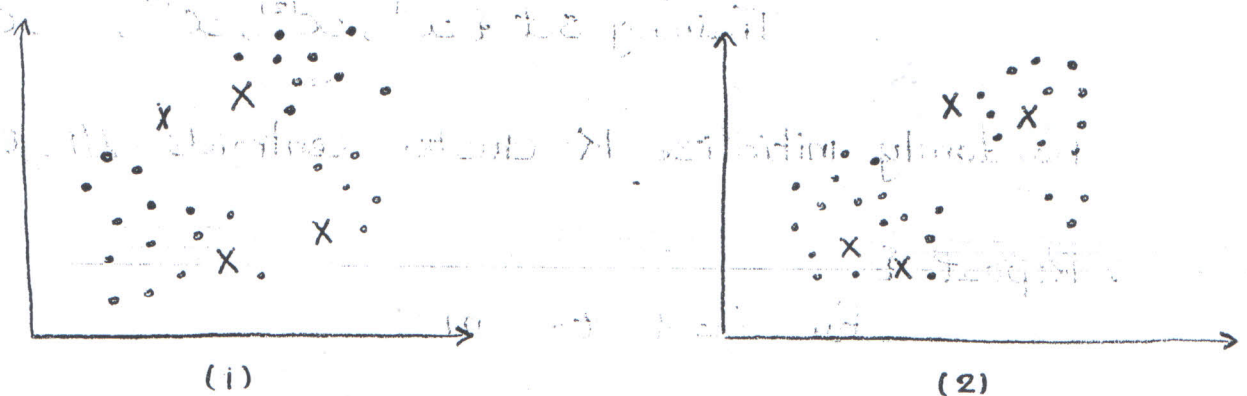
→ We select 2 cluster Centroids in random.



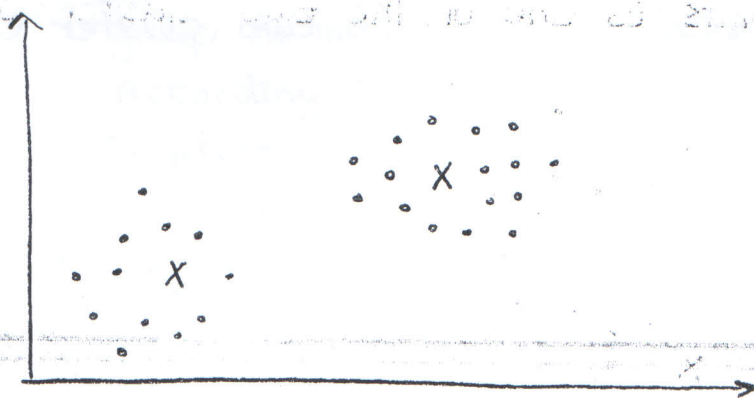
→ In the cluster assignment Step, the algorithm goes through each of the training examples, and depending on whether it's closer to which Cluster Centroid, it is going to assign each of the data points to one of the two cluster Centroids.



→ In the move Centroid Step, the algorithm is going to move the two Cluster Centroids to the average of the points of the same cluster.



→ If you keep running additional iterations of K means from here the cluster Centroids will not change any further. So at this point, K means has converged.



K-means algorithm

- Input :
 - K (number of clusters (groups))
 - Training set $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in$

Repeat &

The cluster assignment step → For $i = 1$ to (m)

↓ number of training sets

Number of cluster which the point will belong to $\leftarrow C^{(i)} := \text{index (From } \overset{\text{no. of clusters}}{1 \text{ to } K}) \text{ of cluster centroid closest to } x^{(i)}$

$$\min_K \|x^{(i)} - \mu_K\|^2$$

→ choose cluster according to min distance between $x^{(i)}$ and cluster centroid

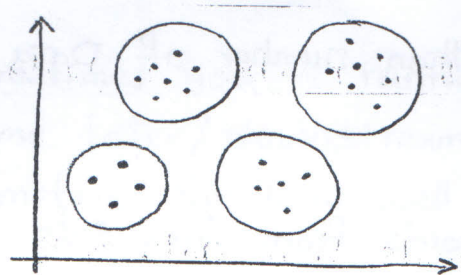
The move centroid step → For $K = 1$ to K

$\mu_K :=$ average (mean) of points assigned to cluster K.

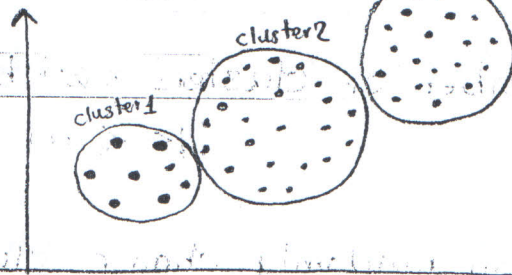
→ Choose new centroid location where the average of sets assigned to this cluster.

}

K-means For non-separated clusters



Separated clusters



non-separated clusters

→ Using K-means algorithm for non-separated clusters, we make sure that each set belongs to the nearest cluster from its location.

K-means Optimization Objective.

Defs

$C^{(i)}$: index of cluster $(1, 2, \dots, K)$ to which example $x^{(i)}$ is currently assigned (cluster number)

μ_k : Cluster Centroid K ($\mu_k \in \mathbb{R}^n$)

$\mu_{C^{(i)}}$: Cluster Centroid of cluster which example $x^{(i)}$ has been assigned. (centroid of cluster that $x^{(i)}$ will go to)

IF $C^{(i)} = 5 \Rightarrow \mu_{C^{(i)}} = \mu_5$

Optimization objective

$$J(c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{C^{(i)}}\|^2$$

Goal: $\min_{c^{(1)}, \dots, c^{(m)}} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

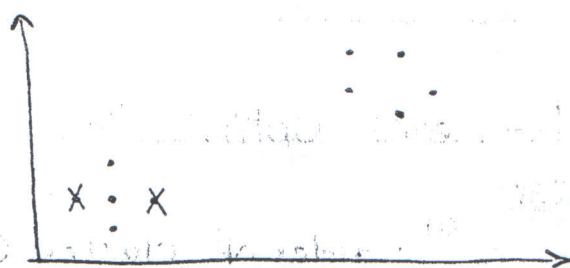
μ_1, \dots, μ_K → This means we must get each set belong to the cluster where the distance between $x^{(i)}$ and μ_k is the smallest

Random initialization

- number of clusters must be smaller than number of Data sets
 $K < m$
- We can randomly choose cluster Centroids from data sets or set them from our side of view.
- random initialization of Cluster Centroids is very important to get right answer or not



we get 2 clusters
(right Answer)



but we can't get 2 clusters
(wrong Answer)

From this Example, you might really guess that K-means can end up converging to different solutions depending on exactly how the clusters were initialized.

For $i = 1$ to 100 &

Randomly initialize K-means.

Run K-means. Get $c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k$

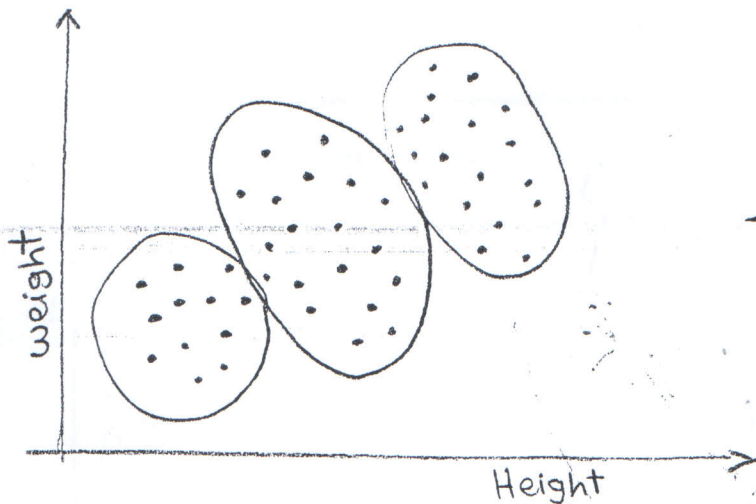
Compute Cost Function

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$$

Pick clustering that gave lowest Cost $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$

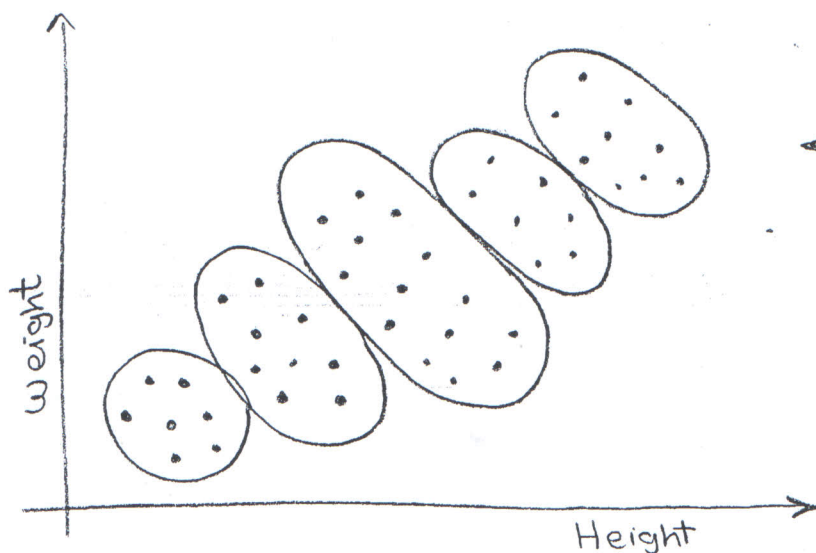
Choosing The Value Of K

- Sometimes, you're running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.



K : Cluster Number

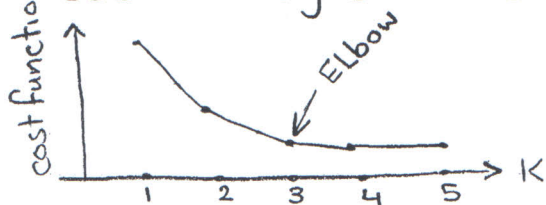
← here, I want to group data into 3 clusters small, mid, large



← here, I want to group data into 5 clusters, small, mid, large, X Large, XX Large

□ ELbow Method : (mathematical method)

determine number of clusters according to Cost Function, as we choose cluster number in which Cost Function change becomes very small. But this method causes waste in time.



we choose $K=3$

1001 1002 1003